

AD-A186 270

ON TWO METHODS OF IDENTIFYING INFLUENTIAL SETS OF
OBSERVATIONS(U) CALIFORNIA UNIV RIVERSIDE DEPT OF
STATISTICS S GHOSH FEB 87 TR-152 AFOSR-TR-87-1244

1/1

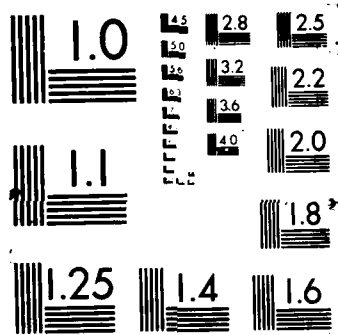
UNCLASSIFIED

AFOSR-87-0048

F/G 12/3

NL





AD-A186 270

DTIC FILE 00

(2)

REPORT DOCUMENTATION PAGE

1a. REPORT SECURITY CLASSIFICATION Unclassified		1b. RESTRICTIVE MARKINGS	
2a. SECURITY CLASSIFICATION AUTHORITY NA		3. DISTRIBUTION/AVAILABILITY OF REPORT Not public release; distribution unlimited.	
2b. DECLASSIFICATION/DOWNGRADING SCHEDULE NA		5. MONITORING ORGANIZATION REPORT NUMBER(S) AFOSR-TR- 87 - 1244	
4. PERFORMING ORGANIZATION REPORT NUMBER(S) Technical Report No. 152		7a. NAME OF MONITORING ORGANIZATION AFOSR/NM	
6a. NAME OF PERFORMING ORGANIZATION University of California, Riverside		7b. ADDRESS (City, State and ZIP Code) Bldg. 410 Bolling AFB DC 20332-6448	
6b. OFFICE SYMBOL (If applicable)		9. PROCUREMENT INSTRUMENT IDENTIFICATION NUMBER AFOSR-87-0048	
6c. ADDRESS (City, State and ZIP Code) Department of Statistics University of California, Riverside Riverside, CA 92521		10. SOURCE OF FUNDING NUMS. PROGRAM ELEMENT NO. 61102F PROJECT NO. 2304 TASK NO. A5 WORK UNIT NO.	
6d. NAME OF FUNDING/SPONSORING ORGANIZATION AFOSR		6e. OFFICE SYMBOL (If applicable) nm	
6e. ADDRESS (City, State and ZIP Code) Bldg. 410 Bolling AFB DC 20332-6448		11. TITLE (Include Security Classification) On Two Methods of Identifying Influential Sets of Observations	
12. PERSONAL AUTHOR(S) Subir Ghosh		13a. TYPE OF REPORT Interim	
13b. TIME COVERED FROM 12/86 TO 02/87		14. DATE OF REPORT (Yr., Mo., Day) February 1987	
15. PAGE COUNT 11		16. SUPPLEMENTARY NOTATION submitted to Statistics and Probability Letters.	

17. COSATI CODES			18. SUBJECT TERMS (Continue on reverse if necessary and identify by block number) Cook's Measure, Design, Influential Observations, Linear Models, Robustness, Unavailable Observations
FIELD	GROUP	SUB. GR.	

19. ABSTRACT (Continue on reverse if necessary and identify by block number)
In this paper two new measures are proposed to identify influential sets of observations at the design state in view of prediction and fitting. A relationship is established between one of proposed measures and the Cook's measure at the inference stage.

DTIC
ELECTE
OCT 15 1987
S D

20. DISTRIBUTION/AVAILABILITY OF ABSTRACT UNCLASSIFIED/UNLIMITED <input checked="" type="checkbox"/> SAME AS RPT. <input type="checkbox"/> DTIC USERS <input type="checkbox"/>		21. ABSTRACT SECURITY CLASSIFICATION Unclassified	
22a. NAME OF RESPONSIBLE INDIVIDUAL Major Brian W. Woodruff		22b. TELEPHONE NUMBER (Include Area Code) (202) 767-5027	
		22c. OFFICE SYMBOL AFOSR/NM	

AFOSR-TR. 87 - 1244

ON TWO METHODS OF IDENTIFYING INFLUENTIAL SETS OF OBSERVATIONS*

by

Subir Ghosh

Department of Statistics
University of California
Riverside, California 92521

In this paper two new measures are proposed to identify influential sets of observations at the design stage in view of prediction and fitting. A relationship is established between one of proposed measures and the Cook's measure at the inference stage.

AMS 1970 Subject Classifications: Primary and Secondary 62J05, 62K15

Short Running Title: Identifying Influential Observations

Keywords and Phrases: Cook's Measure, Design, Influential Observations, Linear Models, Robustness, Unavailable Observations.

* The work of the author is sponsored by the Air Force Office Of Scientific Research under Grant AFOSR-87-0048.

1. Introduction

A set of observations under a design is said to be influential in this paper if the set affects not only the fitting of the model to the data but also the prediction in terms of the fitted model. In the problem of identifying sets of t (a positive integer) influential observations, we assume the underlying design is robust against the unavailability of any t observations [Chosh (1979)]. We first explain this concept by considering the standard linear model

$$E(\underline{y}) = X\underline{\beta} , \quad (1)$$

$$V(\underline{y}) = \sigma^2 I , \quad (2)$$

$$\text{Rank } X = p , \quad (3)$$

where $\underline{y}(N \times 1)$ is a vector of observations, $X(N \times p)$ is a known matrix, $\underline{\beta}(p \times 1)$ is a vector of fixed unknown parameters and σ^2 is a constant which may or may not be known. Let d be the underlying design corresponding to \underline{y} . The design d is assumed to be robust against the unavailability of any t observations in the sense that the parameters in $\underline{\beta}$ are still unbiasedly estimable when any t observations in \underline{y} are unavailable. There are $\binom{N}{t}$ possible sets of t observations. The idea of robustness of designs against unavailability of data is fundamental in measuring the influence of a set of observations.

We first measure the influence of a set of t observations by assuming the observations in the set unavailable and then calculating the sum of variances of their predicted values from the remaining $(N-t)$ observations. The largest value of the sum indicates the corresponding set of t observations is the most influential in terms of precise



<input checked="checked" type="checkbox"/>
<input type="checkbox"/>
<input type="checkbox"/>
<input type="checkbox"/>
<input type="checkbox"/>
<input type="checkbox"/>
<input type="checkbox"/>
<input type="checkbox"/>
<input type="checkbox"/>
<input type="checkbox"/>

A-1

prediction of unavailable observations. We also measure the influence of a set of t observations by assuming them unavailable and then calculating the sum of squares of the elements of the covariance matrix between the least squares fitted values of the remaining $(N-t)$ observations and a complete set of orthonormal linear functions of \underline{y} with zero expectations. The largest value of the sum of squares indicates the corresponding set of t observations is the most influential.

The importance of knowing the influential set of observations at the design stage is that (1) we can assess the influence of a set of unavailable observations in the planned analysis, (2) in the case of deficit of budget during a long term experiment using the robust design where it may be a good idea not to collect observations which are least influential.

2. First Method

We denote the i th set of t observations in \underline{y} by $\underline{y}_2^{(i)}$; and the remaining observations in \underline{y} by $\underline{y}_1^{(i)}$; the corresponding submatrices of X by $X_2^{(i)}$ and $X_1^{(i)}$; the resulting design when t observations in the i th set are unavailable by $d^{(i)}$, $i=1, \dots, \binom{N}{t}$. The least squares estimators of $\underline{\beta}$ under d and $d^{(i)}$ are $\hat{\underline{\beta}}_d = (X'X)^{-1}X'\underline{y}$ and $\hat{\underline{\beta}}_d^{(i)} = (X_1^{(i)'}X_1^{(i)})^{-1}X_1^{(i)'}\underline{y}_1^{(i)}$. We write the fitted values of \underline{y} under d and $d^{(i)}$ as $\hat{\underline{y}}_d = \hat{\underline{\beta}}_d$ and $\hat{\underline{y}}_{d^{(i)}} = X\hat{\underline{\beta}}_d^{(i)}$. When t observations in the i th set are unavailable, the predicted values of unavailable observations $\underline{y}_2^{(i)}$ from available observations are the elements in $\hat{\underline{y}}_2^{(i)} = X_2^{(i)}\hat{\underline{\beta}}_d^{(i)}$. The reliability of these estimators can be judged by $V(\hat{\underline{y}}_2^{(i)}) = \sigma^2 X_2^{(i)}(X_1^{(i)'}X_1^{(i)})^{-1}X_2^{(i)'}.$ The first measure of influence of $\underline{y}_2^{(i)}$ is defined as

$$I_1(y_2^{(i)}) = \text{Trace } V(\hat{y}_2^{(i)}) . \quad (4)$$

The smallest value of $I_1(y_2^{(i)})$, $i=1, \dots, (N)$, for $i=u$, indicates that the u th set of t observations is the least influential in terms of precise prediction of unavailable observations. On the other hand the largest value of $I_1(y_2^{(i)})$, $i=1, \dots, (N)$, for $i=w$, indicates the w th set of t observations is the most influential.

We denote $B_{(i)} = I_t + X_2^{(i)}(X_1^{(i)'} X_1^{(i)})^{-1} X_2^{(i)'}$. It can be checked that

$$B_{(i)}^{-1} = I_t - X_2^{(i)}(X'X)^{-1} X_2^{(i)'}, \quad (5)$$

$$\hat{\beta}_{d(i)} - \hat{\beta}_d = (X'X)^{-1} X_2^{(i)'} B_{(i)} (X_2^{(i)} \hat{\beta}_d - y_2^{(i)}), \quad (6)$$

$$E(X_2^{(i)} \hat{\beta}_d - y_2^{(i)})(X_2^{(i)} \hat{\beta}_d - y_2^{(i)})' = \sigma^2 B_{(i)}^{-1}. \quad (7)$$

We denote the i th observations in y by y_i and the i th row in X by x_i' , $i=1, \dots, N$.

Theorem 1 For any design.

$$\sigma^{-2} I_1(y_2^{(i)}) \geq \sum_{i \in \{i_1, \dots, i_t\}} \frac{x_i'(X'X)^{-1} x_i}{1 - x_i'(X'X)^{-1} x_i}, \quad (8)$$

where the i_1, \dots, i_t rows of X are rows of $X_2^{(i)}$.

Proof. It follows from $B_{(i)}$ and $B_{(i)}^{-1}$ given in (5) that for

$i = i_1, \dots, i_t$,

$$1 + x_i' (X_1^{(i)'} X_1^{(i)})^{-1} x_i \geq \frac{1}{1 - x_i'(X'X)^{-1} x_i},$$

i.e.,

$$x_i' (X_1^{(i)'} X_1^{(i)})^{-1} x_i \geq \frac{x_i'(X'X)^{-1} x_i}{1 - x_i'(X'X)^{-1} x_i}.$$

The rest is easy.

Theorem 2 If for a design, the individual observations are equally influential then

$$I_1(y_1) = \frac{p\sigma^2}{(N-p)} \quad (9)$$

Proof. When the individual observations are equally influential, $I_1(y_1)$ is a constant independent of i for $t=1$ and thus $x_2^{(i)} = x_1$ and $x_1'(x_1^{(i)'}x_1^{(i)})^{-1}x_1$ is a constant independent of i . This in turn implies from (5) that for $t=1$, $x_1'(X'X)^{-1}x_1$ is a constant independent of i . We know $\text{trace } X(X'X)^{-1}X' = p$ and thus $x_1'(X'X)^{-1}x_1 = \frac{p}{N}$. From (5), we get $x_1'(x_1^{(i)'}x_1^{(i)})^{-1}x_1 = \frac{p}{(N-p)}$ and hence the result.

Theorem 3 If for a design, the individual observations are equally influential, then

$$I_1(y_2^{(i)}) \geq \frac{p\sigma^2 t}{(N-p)} \quad (10)$$

Proof. For $t=1$ and equally influential individual observations, $x_1'(X'X)^{-1}x_1 = \frac{p}{N}$ and hence the result in (10) follows from (8). From (9) and (10), we observe that for a design with equally influential individual observations $I_1(y_2^{(i)}) \geq t I_1(y_1)$.

3. Second Method

Let $Z((N-p) \times N)$ be a matrix such that $\text{Rank } Z = (N-p)$, $ZX = 0$ and $ZZ' = I$. It can be seen that $\text{Cov}(\hat{y}_d, z_y) = 0$. This implies that \hat{y}_d has the minimum variance within the class of all unbiased estimators of $E(\hat{y}_d)$ under (1-3). When t observations in the i th set are unavailable, the least squares fitted values are $\hat{y}_1^{(i)} = x_1^{(i)'} \hat{\beta}_{d(i)}$. We denote the submatrices of Z corresponding to $x_1^{(i)}$ and $x_2^{(i)}$ by $Z_1^{(i)}$ and $Z_2^{(i)}$. It follows that $\text{Cov}(\hat{y}_1^{(i)}, z_y) = \sigma^2 [x_1^{(i)}(x_1^{(i)'}x_1^{(i)})^{-1}x_1^{(i)'}Z_1^{(i)'}]$. The

further $\text{Cov}(\hat{y}_1^{(i)}, z_y)$ is away from the null matrix, the more influential is the set of t observations $y_2^{(i)}$. We thus have the second measure of influence as

$$I_2(y_2^{(i)}) = \sigma^{-2} [\text{Sum of squares of elements in } \text{Cov}(\hat{y}_1^{(i)}, z_y)] . \quad (11)$$

We now show some similarities between our two measures of influence

$I_1(y_2^{(i)})$ and $I_2(y_2^{(i)})$.

Theorem 4 The following is true for $i=1, \dots, \binom{N}{t}$,

$$v(z_1^{(i)} \hat{y}_1^{(i)}) = v(z_2^{(i)} \hat{y}_2^{(i)}) .$$

Proof. We observe that $z_1^{(i)} x_1^{(i)} + z_2^{(i)} x_2^{(i)} = 0$ and hence $v(z_1^{(i)} \hat{y}_1^{(i)}) = \sigma^2 z_1^{(i)} x_1^{(i)} (x_1^{(i)'} x_1^{(i)})^{-1} x_1^{(i)'} z_1^{(i)'} = \sigma^2 z_2^{(i)} x_2^{(i)} (x_1^{(i)'} x_1^{(i)})^{-1} x_2^{(i)'} z_2^{(i)'} = v(z_2^{(i)} \hat{y}_2^{(i)})$.

Theorem 5 The following is true.

$$I_2(y_2^{(i)}) = \text{Trace } v(z_2^{(i)} \hat{y}_2^{(i)}) . \quad (13)$$

Proof. It can be seen that

$$\begin{aligned} I_2(y_2^{(i)}) &= \sigma^2 \text{Trace } x_1^{(i)} (x_1^{(i)'} x_1^{(i)})^{-1} x_1^{(i)'} z_1^{(i)'} z_1^{(i)} x_1^{(i)} (x_1^{(i)'} x_1^{(i)})^{-1} x_1^{(i)'} \\ &= \sigma^2 \text{Trace } z_1^{(i)} x_1^{(i)} (x_1^{(i)'} x_1^{(i)})^{-1} x_1^{(i)'} z_1^{(i)'} \\ &= \text{Trace } v(z_1^{(i)} \hat{y}_1^{(i)}) \\ &= \text{Trace } v(z_2^{(i)} \hat{y}_2^{(i)}) \end{aligned}$$

Corollary We have

$$I_2(\hat{y}_2^{(i)}) = \text{Trace } [v(\hat{y}_2^{(i)})][z_2^{(i)'} z_2^{(i)}] . \quad (14)$$

The equation (13) displays the similarity between two measures of

influence $I_1(y_2^{(i)})$ and $I_2(y_2^{(i)})$. Although the matrix Z is not unique, it can be checked that $I_2(y_2^{(i)})$ is unique for all choices of the matrix Z .

4. Relationship

Cook (1977) proposed a distance function between \hat{y}_d and $\hat{y}_{d(i)}$, popular as Cook's distance, at the inference stage as

$$D_i = \frac{(\hat{y}_{d(i)} - \hat{y}_d)'(\hat{y}_{d(i)} - \hat{y}_d)}{ps_d^2}, \quad (15)$$

where $(N-p)s_d^2 = (\underline{y} - \hat{\underline{y}}_d)'(\underline{y} - \hat{\underline{y}}_d)$. The Cook's distance D_i measures the degree of influence of t observations in the i th set on the estimation of $\underline{\beta}$. We now show that our first measure of influence $I_1(\underline{y}_2^{(i)})$ is in fact related to D_i .

Theorem 6 From (4) and (15), we have

$$E(ps_d^2 D_i) = I_1(\underline{y}_2^{(i)}). \quad (16)$$

Proof. We get from (6)

$$(\hat{\underline{y}}_{d(i)} - \hat{\underline{y}}_d)'(\hat{\underline{y}}_{d(i)} - \hat{\underline{y}}_d) = (\underline{x}_2^{(i)} \hat{\underline{\beta}} - \underline{y}_2^{(i)})'(B_{(i)}^2 - B_{(i)})(\underline{x}_2^{(i)} \hat{\underline{\beta}} - \underline{y}_2^{(i)}).$$

It now follows from (7), (15) and (17) that

$$\begin{aligned} E(ps_d^2 D_i) &= \sigma^2 \text{Trace} (B_{(i)} - J_t) \\ &= \sigma^2 \text{Trace} \underline{x}_2^{(i)} (\underline{x}_1^{(i)'} \underline{x}_1^{(i)})^{-1} \underline{x}_2^{(i)'} \\ &= I_1(\underline{y}_2^{(i)}). \end{aligned}$$

This completes the proof.

5. Examples

Consider a 2^4 factorial experiment in a completely randomised set up and suppose the elements in $\underline{\beta}$ are the general mean, the main effects and the 2-factor interactions. The 3-factor and higher order interactions are assumed to be zero. Thus $p=11$. The treatments are denoted by $(x_1 x_2 x_3 x_4)$, $x_i=0,1, i=1,2,3,4$. For brevity, we indicate a treatment by the

positions where the level 1 is occurring. For example, the treatment (1100) is denoted by 12. The treatment (0000) is denoted by 0.

Design I

Consider a design with 15 treatments and we write the treatments in the order (0,1,2,3,4,12,13,14,23,24,34,123,124,134,234). Note that the elements in y , the rows of X and the columns of Z correspond to this ordering. The matrix Z is given below

$$Z = (.25) \begin{bmatrix} 4a & -3a & -3a & -3a & -3a & 2a & 2a & 2a & 2a & 2a & 2a & -a & -a & -a & -a \\ 0 & 1 & -1 & 1 & -1 & 0 & -2 & 0 & 0 & 2 & 0 & 1 & -1 & 1 & -1 \\ 0 & 1 & -1 & -1 & 1 & 0 & 0 & -2 & 2 & 0 & 0 & -1 & 1 & 1 & -1 \\ 0 & 1 & 1 & -1 & -1 & -2 & 0 & 0 & 0 & 0 & 2 & 1 & 1 & -1 & -1 \end{bmatrix},$$

where $a = (1/\sqrt{5})$. The design is robust against the unavailability of any two observations [Ghosh (1979)].

Table I
Influences of Individual Observations Under Design I

Observations	$\sigma^{-2}I_1$	$\sigma^{-2}I_2$
y_1	4.000	.800
$y_i, i=2, \dots, 11$	2.333	.700
$y_i, i=12, \dots, 15$	4.000	.800

Table II

Influences of Pairs of Observations Under Design I

Observations	$\sigma^{-2}_{I_1}$	$\sigma^{-2}_{I_2}$
$(y_1, y_1), i=2, \dots, 15; (y_2, y_{15});$ $(y_3, y_{14}); (y_5, y_{12}); (y_6, y_1) i=2, 13;$ $(y_7, y_1) i=12, 14; (y_8, y_1) i=13, 14;$ $(y_9, y_1) i=12, 15; (y_{10}, y_1) i=13, 15;$ $(y_{11}, y_1) i=14, 15.$	11.333	1.500
$(y_1, y_1) i=6, \dots, 11; (y_2, y_1) i=12, 13, 14;$ $(y_3, y_1) i=12, 13, 15; (y_4, y_1) i=12, 14, 15;$ $(y_5, y_1) i=13, 14, 15; (y_6, y_1) i=14, 15;$ $(y_7, y_1) i=13, 15; (y_8, y_1) i=12, 15;$ $(y_9, y_1) i=13, 14; (y_{10}, y_{14});$ $(y_{11}, y_1) i=12, 13$	8.000	1.500
$(y_1, y_1) i=12, \dots, 15;$ $(y_{12}, y_1) i=13, 14, 15;$ $(y_{13}, y_1) i=14, 15;$ (y_{14}, y_{15})	8.667	1.600

Table II (Continued)

Influences of Pairs of Observations Under Design I

Observations	$\sigma^{-2}_{I_1}$	$\sigma^{-2}_{I_2}$
$(y_2, y_1)_{i=3,4,5,9,10,11};$ $(y_3, y_1)_{i=4,5,7,8,11};$ $(y_4, y_1)_{i=5,6,8,10};$ $(y_5, y_1)_{i=6,7,9};$ $(y_6, y_1)_{i=7,8,9,10};$ $(y_7, y_1)_{i=8,9,11}; (y_8, y_1)_{i=10,11};$ $(y_9, y_1)_{i=10,11}; (y_{10}, y_{11}).$	4.857	1.400
$(y_2, y_1)_{i=6,7,8}; (y_3, y_1)_{i=6,9,10};$ $(y_4, y_1)_{i=7,9,11};$ $(y_5, y_1)_{i=8,10,11};$ $(y_6, y_{11}), (y_7, y_{10});$ $(y_8, y_9), (y_{10}, y_{12})$	10.000	1.400

We find that under Design I, any of $y_i, i=2, \dots, 11$ is the least influential w.r.t. both I_1 and I_2 . Any pair of observations with $\sigma^{-2}_{I_1}$ equals 11.333 is the most influential w.r.t. I_1 . On the other hand, any pair of observations with $\sigma^{-2}_{I_2}$ equals 1.600 is the most influential w.r.t. I_2 . The variability in values of I_2 is so small that it is very hard to assess the influence w.r.t. I_2 under this design.

Design II

We consider a complete 2^4 factorial experiment with treatments written in the order (1234, 0, 1,2,3,4,12,13,14,23,24,34,123,124,134, 234). The matrix Z is as follow

$$Z = (.25) \begin{bmatrix} 1 & 1 & -1 & -1 & -1 & -1 & 1 & 1 & 1 & 1 & 1 & 1 & -1 & -1 & -1 & -1 \\ 0 & 0 & 1 & -1 & 1 & -1 & 0 & -2 & 0 & 0 & 2 & 0 & 1 & -1 & 1 & -1 \\ 0 & 0 & 1 & -1 & -1 & 1 & 0 & 0 & -2 & 2 & 0 & 0 & -1 & 1 & 1 & -1 \\ 0 & 0 & 1 & 1 & -1 & -1 & -2 & 0 & 0 & 0 & 0 & 2 & 1 & 1 & -1 & -1 \\ 2 & -2 & 1 & 1 & 1 & 1 & 0 & 0 & 0 & 0 & 0 & 0 & -1 & -1 & -1 & -1 \end{bmatrix},$$

This design is robust against the unavailability of any three observations [Ghosh (1979)]. It can be checked that $I_1(y_i) = 2.2\sigma^2$ and $I_2(y_i) = (.6875)\sigma^2$ for $i=1, \dots, 16$. Thus the individual observations are equally influential w.r.t. both I_1 and I_2 . For any pair of observations corresponding to treatments with zero or three levels in common, the value of $\sigma^{-2}I_1$ is 8.000. For every other pair of observations, the value of $\sigma^{-2}I_1$ is 4.667. We therefore see the validity of the equation (10) in Theorem 3 for this example since $2 I_1(y_i) = 4.4\sigma^2$. The value of $\sigma^{-2}I_2$ for any pair of observations is a constant 1.375. The remark on I_2 for Design I also holds for Design II.

REFERENCES

- Chatterjee, S. and Hadi, A. S. (1986). Influential Observations, high leverage points and outliers in linear regression. Statistical Science, 3, 379-416
- Cook, R. D. (1977). Detection of influential observations in linear regression. Technometrics, 22, 495-508.
- Cook, R. D. and Weisberg, S. (1982). Residuals and influence in regression. Technometrics, 23, 21-26.
- Draper, N. R. and John, J. A. (1981). Influential observations and outliers in regression. Technometrics, 23, 21-26.
- Ghosh, S. (1979). On robustness of designs against incomplete data, Sankhyā, B, Pts 3 and 4, 204-208.
- Kiefer, J. (1959). "Optimum experimental designs." Roy. Statist. Soc. Ser. B, 21, 273-319.

END

12-87

DTIC